

An annotated chromosome-scale reference genome for Eastern black-eared wheatear (*Oenanthe melanoleuca*)

Valentina Peona^{1,*}, Octavio Manuel Palacios-Gimenez^{1,2,3*}, Dave Lutgen^{4,2,5,*}, Remi André Olsen⁶, Niloofar Alaei Kakhki², Pavlos Andriopoulos⁷, Vasileios Bontzorlos⁸, Manuel Schweizer^{9,4}, Alexander Suh^{1,10}, Reto Burri^{5,4,3}

¹ Department of Organismal Biology – Systematic Biology, Science for Life Laboratory, Evolutionary Biology Centre, Uppsala University, 75236 Uppsala, Sweden

² Department of Population Ecology, Institute of Ecology and Evolution, Friedrich Schiller University Jena, 07743 Jena, Germany

³ German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, 04103 Leipzig, Germany

⁴ Institute of Ecology and Evolution, University of Bern, 3012 Bern, Switzerland

⁵ Swiss Ornithological Institute, CH-6204 Sempach, Switzerland

⁶ Science for Life Laboratory, Department of Biochemistry and Biophysics, Stockholm University, 17165 Solna, Sweden

⁷ Section of Ecology and Systematics, Department of Biology, National and Kapodistrian University of Athens, 15772 Athens, Greece

⁸ TYTO – Association for the Management and Conservation of Biodiversity in Agricultural Ecosystems, 41335 Larisa, Greece

⁹ Natural History Museum Bern, 3005 Bern, Switzerland

¹⁰ School of Biological Sciences, University of East Anglia, NR4 7TU Norwich, United Kingdom

* These authors contributed equally to the work.

Keywords

Birds, open-habitat chats, *Oenanthe melanoleuca*, *Oenanthe hispanica*-complex, transcriptome, repeat content, transposable elements

Running head

Eastern black-eared wheatear genome

Correspondence

Reto Burri, Swiss Ornithological Institute, Seerose 1, CH-6204 Sempach, Switzerland

Email: reto.burri@vogelwarte.ch

ORCID

Valentina Peona: 0000-0001-5119-1837

Octavio Manuel Palacios-Gimenez: 0000-0002-1472-9949

Dave Lutgen: 0000-0003-0793-3930

Pavlos Andriopoulos: 0000-0002-5377-2974

Vasileios Bontzorlos: 0000-0002-1276-3385

Manuel Schweizer: 0000-0002-7555-8450

1 Alexander Suh: 0000-0002-8979-9992

2 Reto Burri: 0000-0002-1813-0079

3

4 **Abstract**

5 Pervasive convergent evolution and in part high incidences of hybridization distinguish wheatears
6 (songbirds of the genus *Oenanthe*) as a versatile system to address questions at the forefront of
7 research on the molecular bases of phenotypic and species diversification. To prepare the genomic
8 resources for this venture, we here generated and annotated a chromosome-scale assembly of the
9 Eastern black-eared wheatear (*O. melanoleuca*). This species is part of the *O. hispanica*-complex that
10 is characterized by convergent evolution of plumage coloration and high rates of hybridization. The
11 long-read-based male nuclear genome assembly comprises 1.04 Gb in 32 autosomes, the Z
12 chromosome, and the mitogenome. The assembly is highly contiguous (contig N50: 12.6 Mb; scaffold
13 N50: 70 Mb), with 96 % of the genome assembled at chromosome level and 95.5 % BUSCO
14 completeness. The nuclear genome was annotated with 18,143 protein-coding genes and 31,333
15 mRNAs (annotation BUSCO completeness: 98.0 %), and about 10 % of the genome consists of
16 repetitive DNA. The annotated chromosome-scale reference genome of Eastern black-eared wheatear
17 provides a crucial resource for research into the genomics of adaptation and speciation in an
18 intriguing group of passerines.

19

20 **Introduction**

21 Wheatears of the genus *Oenanthe* and their relatives – together referred to as “open-habitat chats” –
22 are a group of songbirds that display several remarkable characteristics distinguishing them as a
23 versatile system to address key questions on the evolution of phenotypes and formation of species.
24 Many phenotypes, including multiple conspicuous colour ornaments, seasonal migration, and sexual
25 dimorphism appear independently in multiple branches within open-habitat chats, suggesting a high
26 incidence of convergent evolution (Alaei Kakhki et al. in press; Aliabadian et al. 2012; Schweizer et al.
27 2019). Furthermore, hybridization is observed in several species complexes and occurs at notably
28 high rates in the *O. hispanica*-complex that consists of four currently recognized taxa (Schweizer et al.
29 2019): Western black-eared wheatear (*O. hispanica*), pied wheatear (*O. pleschanka*), cyprus wheatear
30 (*O. cyprica*), and Eastern black-eared wheatear (*O. melanoleuca*; **Fig. 1**). Pied and Eastern black-
31 eared wheatear hybridize pervasively at the western shores of the Black Sea, in the Caucasus, and in
32 the Alborz mountains of northern Iran (Haffer 1977; Panov 2005). The resulting introgression
33 reaches beyond the hybrid zones (Schweizer et al. 2019), and hybrid zones themselves sport admixed
34 phenotypes that display combinations of plumage colour phenotypes divergent between species
35 (mantle and neck-side coloration) (Haffer 1977; Panov 2005). Finally, a phenotype divergently
36 expressed between many wheatear species, black-or-white throat coloration, segregates as
37 polymorphisms in three species of the *O. hispanica*-complex. Once a high-quality reference genome is
38 available, this polymorphism and the recombination of mantle and neck-side coloration in hybrids
39 provide an excellent opportunity to map these phenotypes to the genome (Buerkle and Lexer 2008)
40 and study their convergent evolution across open-habitat chats. Furthermore, hybridization in

1 several geographic regions enables insights into common or idiosyncratic patterns of evolution under
2 hybridization (Gompert et al. 2017).

3 Here, we describe the *de novo* assembly and annotation of a chromosome-scale reference
4 genome for the Eastern black-eared wheatear (*O. melanoleuca*). The assembly includes models for 32
5 autosomes, the Z chromosome and the mitogenome that together cover 90 % of the k-mer-based
6 genome size estimate (94 % with unplaced scaffolds included); it is highly contiguous with a scaffold
7 N50 of 70 Mb and BUSCO completeness score of 95.5 %. This reference genome enables genomic
8 research into the evolutionary history of phenotypic and species diversification in wheatears and
9 their close relatives.

10 **Material and Methods**

11 **Sampling, tissue preservation, and nucleic acid extraction**

12 To obtain optimal starting material for a reference individual, we freshly sampled a male Eastern
13 black-eared wheatear (*Oenanthe melanoleuca*) well outside known hybrid zones (Haffer 1977; Panov
14 2005) in Galaxidi, Greece (sampling permit no. 181968/989, issued by the Ministry of Environment
15 and Energy, General Secretariat of Environment, General Directorate of Forests and Forest
16 Environment, Directorate of Forest Management, Department of Wildlife and Game Management;
17 export permit no. 55980/1575, Regional CITES management authority Attika). For this purpose, we
18 sampled about 100 µl of blood from the brachial vein, and, after euthanizing the bird, we extracted all
19 tissues possible. Tissues were immediately snap-frozen in liquid nitrogen. Throughout transportation
20 and storage preceding DNA extraction, the samples were kept at a temperature below -80° C.

21 To obtain ultra-high molecular weight (UHMW) DNA from the reference individual, NGI
22 Uppsala (Sweden) extracted DNA from the blood sample using the Bionano Prep™ Blood and Cell
23 Culture DNA Isolation Kit (Bionano, San Diego, USA). Electrophoresis on a Femto Pulse instrument
24 showed a mean DNA fragment length of about 200 kb, with fragments reaching up to 800 kb.

25 To prepare muscle tissue for Hi-C sequencing library preparation, we pulverized breast
26 muscle tissue from the reference individual in a mortar. To avoid unfreezing of the tissue powder, the
27 procedure was carried out in a climate chamber at 4°C under regular addition of liquid nitrogen.

28 To prepare RNA for full-length transcript sequencing, we extracted total RNA from eight snap-
29 frozen tissues kept at -80°C (brain, breast muscle, heart, kidney, liver, lung, spleen, and testis) using
30 the RNeasy Mini Kit (Qiagen; Hombrechtikon, Switzerland) according to the manufacturer's
31 instructions. RNA quality was assessed with a Fragment Analyzer (Agilent). RNA from spleen showed
32 considerable degradation and was excluded from further analyses.

33 ***De novo* genome sequencing, and reference genome assembly and annotation**

34 **Assembly strategy and data acquisition**

35 To obtain a chromosome-scale reference genome, our strategy largely followed the multiplatform
36 approach recommended by Peona et al. (2021). In brief, it consisted of (i) a phased primary assembly
37 based on long reads (ii) polishing and scaffolding of the primary assembly with linked-read

1 sequencing data, and (iii) scaffolding of the secondary assembly with proximity ligation (Hi-C)
2 information.

3 To this end, we obtained a total of 215 Gb (unique coverage 151 Gb) Pacific Biosciences
4 (PacBio) long-read sequence data, 54 Gb linked-read sequence data, and 83 Gb Hi-C data. NGI Uppsala
5 (Sweden) prepared a PacBio library from UHMW DNA using the SMRTbell Template Prep Kit 1.0 and
6 sequenced this library on 18 SMRT Cells 1M v3 on a PacBio Sequel instrument (Sequel Binding Kit
7 3.0, Sequel Sequencing Plate 3.0). PacBio long-read data was initially processed using SMRT Link v6.
8 A linked-read sequencing library was prepared using the 10X Genomics Chromium Genomic Kit (from
9 the same DNA extraction as used for PacBio sequencing; 10X Genomics, Inc., Pleasanton, CA, USA; Cat
10 No. 120215), and a Hi-C library was prepared with the Dovetail Omni-C kit (Scotts Valley, CA, USA;
11 Cat No. 21005). The linked-read and Hi-C libraries were prepared and sequenced on a NovaSeq 6000
12 instrument (S4 lane, 150 bp paired-end reads) at the facilities of NGI Stockholm (Sweden).

13 **Genome size estimation**

14 We estimated genome size by counting k-mer frequency of the quality-checked 10X Genomics linked
15 reads. To this end, we first trimmed 22 bp from all 10X Genomics linked reads using fastp (Chen et al.
16 2018) to remove indices from R1 reads and keep symmetric read lengths for the R2 reads. We then
17 counted k-mers of size 21 using jellyfish 2.2.10 (Marçais and Kingsford 2011) and used GenomeScope
18 (Vurture et al. 2017) to estimate genome size from k-mer count histograms.

19 ***De novo* genome assembly**

20 We assembled the PacBio long reads into the phased primary assembly using the Falcon Unzip 0.5
21 assembler (Chin et al. 2016), followed by polishing with Arrow 1.9.0. Before assembly polishing, we
22 masked repeat regions of the phased primary assembly with RepeatMasker 4.1.0 (Smit et al. 1996-
23 2010) using a custom repeat library (Boman et al. 2019; Peona et al. 2021; Suh et al. 2018;
24 Weissensteiner et al. 2020) to make accurate assembly corrections without overcorrecting large
25 repeats. We then polished the masked assembly with two rounds in Pilon v1.22 (Walker et al. 2014)
26 with the parameter "--fix indels" using the reference individual's linked-read data. To purge duplicate
27 scaffolds from the assembly, we ran purge_dups 1.2.5 (Guan et al. 2020) on the polished assembly.
28 Prior to scaffolding with linked-read data, we split potential mis-assemblies with reference-individual
29 linked-read data using Tigrint 1.2.4 (Jackman et al. 2018). With the aim to scaffold the polished
30 remaining contigs, we applied ARCS 1.2.2 and LINKS 2.0.0 using the reference individual's linked-read
31 data using default parameters (Warren et al. 2015; Yeo et al. 2018).

32 To further scaffold the assembly, we applied the 3D-DNA pipeline (Dudchenko et al. 2017) to
33 join the sequences into chromosomes. We first used Juicer v.1.6 (Durand et al. 2016) to map Hi-C data
34 against the contigs and to filter reads, and then ran the asm-pipeline v.180922 to generate a draft
35 scaffolding.

36 Finally, we corrected mis-assemblies based on the visual inspection of the proximity map
37 using Juicebox 2.13.06 (Robinson et al. 2018). The final chromosome-level assembly was polished
38 with two additional rounds in Pilon as described above.

39 To assess homology of the assembled scaffolds with bird chromosomes, we aligned the final
40 genome assembly to the genomes of collared flycatcher (*Ficedula albicollis*) (FicAlb1.5) (Kawakami

1 et al. 2014), zebra finch (taeGut3.2.4) (Warren et al. 2010), and chicken (GRCg6a) (Bellott et al. 2017)
2 using D-Genies (Cabanettes and Klopp 2018). Chromosomes were named according to homology with
3 these three genomes. In cases, such as chicken chromosomes 1 and 4 that are split to multiple
4 chromosomes in songbirds, the nomenclature in the wheatear genome was adapted to the species
5 whose homologous chromosome matched closest.

6 **Mitogenome assembly**

7 To assemble the mitochondrial genome, we used the MitoFinder 1.4 (Allio et al. 2020) and mitoVGP
8 2.2 (Formenti et al. 2021) pipelines with the published *Oenanthe isabellina* mitochondrial genome
9 (Genbank Accession Number: NC_040290.1) as reference. We ran MitoFinder with the reference
10 individual's short-read data (linked-read data but without making use of the linked-read haplotype
11 information), and with mitoVGP we made joint use of the linked-read and long-read data. From
12 MitoFinder we extracted the longest contig containing all 13 protein coding genes, two rRNA genes
13 and 22 tRNAs annotated by MitoFinder as mitogenome assembly. We annotated both assemblies
14 using the MITOS WebServer (<http://mitos2.bioinf.uni-leipzig.de/index.py>).

15 We then aligned both resulting assemblies with the mitogenomes of isabelline wheatear (*O.*
16 *isabellina*, NC_040290.1) and northern wheatear (*O. Oenanthe*, MN356231.1) using MUSCLE (Edgar
17 2004) in MEGA X (Stecher et al. 2020) and generated a circular mitogenome map using CGView
18 (Stothard and Wishart 2005).

19 **Assembly quality evaluation**

20 To evaluate assembly quality at each assembly step, we estimated basic assembly statistics using
21 QUAST 5.0.2 (Gurevich et al. 2013) and evaluated the completeness of expected gene content in the
22 assembly based on benchmarking universal single-copy orthologs (BUSCO) (Simão et al. 2015) with
23 the avian dataset aves_odb10 (8,338 BUSCO) in BUSCO 5.0.0.

24

25 **Repeat annotation**

26 The final version of the genome assembly was used to *de novo* characterize both interspersed and
27 tandem repeats. For interspersed repeats, we used RepeatModeler2 (Flynn et al. 2020) with the
28 option “-LTR_struct” to obtain an improved characterisation of LTR retrotransposons which are
29 commonly found in avian genomes (Boman et al. 2019; Kapusta and Suh 2017; Peona et al. 2021).
30 The resulting library of raw consensus sequences was filtered from consensus sequences of tandem
31 repeats (for which we ran a specific analysis; see below) and from protein-coding genes using the
32 Snakemake pipeline repeatlib_filtering_workflow v0.1.0 ([https://github.com/
33 NBISweden/repeatlib_filtering_workflow](https://github.com/NBISweden/repeatlib_filtering_workflow)).

34 For tandem repeats, we used RepeatExplorer2 (Novák et al. 2020) to search for satellite DNA
35 (satDNA) sequences using the reference individual's 10X Genomics linked reads. Prior to
36 RepeatExplorer2 graph-based clustering analysis, sequencing reads were pre-processed and checked
37 by quality with FastQC (Babraham Bioinformatics: Cambridge 2012) using the public online platform
38 at <https://repeatexplorer.elixir-cerit-sc.cz>. We processed the reads with the "quality trimming tool",
39 "FASTQ interlacer on the paired end reads", "FASTQ to FASTQ converter", followed by

1 "RepeatExplorer2 clustering" with default parameters. Each reference sequence assembled by
2 RepeatExplorer2 consisted of a monomer of the satDNA consensus sequence. The relative genomic
3 abundance and nucleotide divergence (Kimura-2-parameter distance) of each detected satDNA were
4 estimated by sampling four million read pairs and aligning them to the satDNA library with
5 RepeatMasker 4.1.0 (Smit et al. 1996-2010). The sampled reads were mapped to dimers of satDNA
6 consensus sequences, and for smaller satDNAs, several monomers were concatenated until reaching
7 roughly 150 bp array length. The resulting RepeatMasker *.align* file was then parsed to the script
8 *calcDivergenceFromAlign.pl* from RepeatMasker utils. The relative abundance of each satDNA
9 sequence was then estimated as the proportion of nucleotides aligned with the reference sequence
10 with respect to the total Illumina library size.

11 The RepeatModeler2 library was then merged with the satDNA library produced here and
12 with known avian consensus sequences of transposable elements from Repbase (Bao et al. 2015),
13 Dfam (Storer et al. 2021, 2021), flycatcher, blue-capped cordon-bleu, hooded crow, and paradise crow
14 (Boman et al. 2019; Peona et al. 2021; Suh et al. 2018; Weissensteiner et al. 2020). This library was
15 then used to annotate the genome assembly with RepeatMasker (Smit et al. 1996-2010). The
16 annotation produced was processed with the script *calcDivergenceFromAlign.pl* from RepeatMasker
17 utils to calculate the divergence between repeats and their consensus sequences using the Kimura 2-
18 parameter distance corrected for the presence of CpG sites.

19 **Full-length transcript sequencing and genome annotation**

20 We aimed to establish a high-quality genome annotation based on full-length transcripts. To this end,
21 for each of the abovementioned seven tissues, the NGS platform of the University of Berne,
22 Switzerland, prepared an Iso-Seq library using the SMRTbell Express Template Prep Kit 2.0 (Pacific
23 Biosciences). These seven libraries were then sequenced on three separate SMRT cells 8M,
24 sequencing twice five tissues (brain and testis, lung, muscle, and heart) and once two tissues (liver
25 and kidney) per SMRT cell. Sequencing of these SMRT cells was conducted on a Pacific Biosciences
26 Sequel II instrument at the Genomic Technologies Facility in Lausanne, Switzerland. As the libraries
27 underloaded, five libraries (all but liver and kidney) were jointly sequenced on an additional SMRT
28 cell 8M on a Pacific Biosciences Sequel IIe at the NGS platform of the University of Berne.

29 Circular consensus sequences (CCS), full-length non-chimeric transcripts, and polished high-
30 and low-quality transcripts were obtained by the NGS platform at the University of Bern separately
31 for each run using the Isoseq 3 pipeline (ICS v10.1). Polished full-length isoforms for each sequencing
32 run were merged by tissue and then separately mapped to the reference genome using Minimap v2.2
33 (-ax splice) (Li 2018, 2021). Transcriptome annotations were generated by first collapsing redundant
34 transcripts using TAMA collapse (-x no_cap), before generating open reading frame (ORF) and
35 nonsense-mediated mRNA decay (NMD) predictions using the scripts implemented in TAMA-GO (Kuo
36 et al. 2020) for each of the seven tissues. We then evaluated tissue-specific transcriptome
37 completeness using BUSCO (Simão et al. 2015) with the avian dataset aves_odb10 (8'338 BUSCO) in
38 BUSCO 5.0.0. Additional transcriptome annotation statistics were obtained using the
39 *agat_sp_statistics.pl* script implemented in the AGAT toolkit (Dainat 2019).

40 We annotated the repeat soft-masked genome using GeMoMa 1.9 (Keilwagen et al. 2018;
41 Keilwagen et al. 2019), a homology-based gene prediction tool. This tool is based on the annotation

1 of protein-coding genes and intron position conservation in a reference genome to predict the
 2 annotation of protein-coding genes in the target genome. We used the genomes of chicken
 3 (GCA_016699485.1; International Chicken Genome Sequencing Consortium 2004), zebra finch
 4 (GCA_003957565.2; Warren et al. 2010), silvereye (GCA_001281735.1; Cornetti et al. 2015), and
 5 collared flycatcher (GCA_000247815.2; Ellegren et al. 2012; Kawakami et al. 2014) as references for
 6 the homology-based gene prediction, along with the reference individual's transcriptome obtained
 7 from Iso-Seq data to incorporate RNA evidence for the splice prediction. Using the Extract RNA-seq
 8 Evidence tool implemented in GeMoMa, we obtained intron position and coverage. This information
 9 was fed into the GeMoMa pipeline (GeMoMa.m=200000, AnnotationFinalizer.r=SIMPLE, pc=true, and
 10 o=true) to obtain predicted protein-coding gene models. To account for redundancies/duplicates
 11 resulting from the predicted protein-coding genes potentially stemming from each of the four
 12 reference species, genome annotation completeness was assessed by recomputing BUSCO using the
 13 BUSCOrecomputer tool in GeMoMa.

14 Functional annotation of protein-coding genes was obtained with InterProScan 5.59 (Jones et
 15 al. 2014; Paysan-Lafosse et al. 2022). InterProScan ran with the following settings: *-goterms -*
 16 *iprlookup -appl CDD, COILS, Gene3D, HAMAP, MobiDBLite, PANTHER, Pfam, PIRSF, PRINTS,*
 17 *PROSITEPATTERNS, PROSITEPROFILES, SFLD, SMART, SUPERFAMILY, TIGRFAM*). Predicted protein-
 18 coding genes were further annotated through a protein Blast search (-evaluate 0.000001, -seg yes, -
 19 soft_masking true, -lcase_masking) against the Swiss-Prot database (Uniprot Consortium 2019). We
 20 then merged the predicted protein-coding gene models and the functional annotation using the
 21 *agat_sp_manage_functional_annotation.pl* script, obtained summary statistics using
 22 *agat_sp_statistics.pl* and *agat_sp_functional_statistics.pl*, both implemented in the AGAT toolkit. Gene
 23 ontology (GO-terms) were visualised with WEGO 2.0 (wego.genomics.cn).

24 Results and Discussion

25 Nuclear genome assembly

26 The polished, unzipped primary assembly contained a total of 1,681 contigs, of which all were >25 kb
 27 long and 1,610 were >50 kb long (**Tab. 1**). Total assembly length was 1.29 Gb, with the longest contig
 28 spanning 45.3 Mb, contig N50 of 8.6 Mb, and half of the assembly placed in 35 contigs. Avian BUSCO
 29 were 96.9 % complete, with 90.6 % being single-copy genes (**Tab. 1**).

30 Purging duplicated contigs resulted in an assembly constituted of 381 contigs with a total
 31 assembly length of 1.04 Gb, contig N50 of 13.5 Mb and half of the assembly placed in 23 contigs (**Tab.**
 32 **1**). After this step, BUSCO completeness remained at 96.4 %, but an improvement to nearly 96 %
 33 single-copy BUSCOs was achieved (**Tab. 1**).

34 Starting from an already highly contiguous assembly, the linked-read data did not yield any
 35 scaffolding improvement. Still, Tigmint detected several supposed mis-assemblies and split the
 36 assembly into 451 scaffolds. However, an alignment of the original contigs in D-Genies (Cabanettes
 37 and Klopp 2018) showed that all but one of the original contigs (see below) were collinear with the
 38 collared flycatcher genome. Given this result and that the proximity ligation data would correct mis-
 39 assemblies in subsequent steps, we decided to keep the original contigs except for one aligning to
 40 flycatcher chromosomes 2 and 3. For the latter contig, we used the output of Tigmint that split the

1 contig in line with the alignment. The two split parts covered all but 12,527 bp of the original contig.
2 Visual inspection of the missing sequence showed that it almost entirely consisted of repeats. We left
3 this sequence in the assembly as a separate contig.

4 The proximity ligation information obtained through Hi-C scaffolding corrected a number of
5 scaffolds, resulting in a higher number of scaffolds (588) than the number of contigs it started from
6 (383). However, the scaffolding yielded a highly contiguous chromosome-scale assembly (N50, 69.6
7 Mb; L50, 6) with BUSCO completeness of still >96 % and almost all BUSCOs in single copy (**Tab. 1**).
8 This final assembly contained all macrochromosomes and the majority of microchromosomes usually
9 found in the latest generation of avian genome assemblies (Kapusta et al. 2017; Peona 2021; Rhie et
10 al. 2021). 96 % of the assembly was placed into chromosome models, and the chromosome-only
11 assembly covered still 95.5 % of BUSCO (**Tab. 1**).

12 The final assembly length closely matched the one of previous linked-read-based assemblies
13 of the same species and closely related ones (Lutgen et al. 2020; Schweizer et al. 2019). The genome
14 size estimated from the k-mer distribution of linked reads sequence was between 1.105 and 1.106
15 Gb, with 0.925-0.926 Gb of unique and 0.179-0.180 Gb (16 %) repeat sequence and 0.75-0.76 %
16 heterozygosity (GenomeScope model fit 98-99 %). The full final reference genome assembly thus
17 covered 94 % of the genome size estimate, with 90 % of the estimated genome size placed in
18 chromosomes. 96 % of the assembly were placed in 33 chromosomes with homologs in collared
19 flycatcher, zebra finch and chicken, according to which we adapted the chromosome nomenclature.
20 The differences in genome size estimates based on the k-mer approach and the genome assembly
21 length is likely the result of highly repetitive sequences (e.g., centromeres, telomeres, satDNAs) that
22 collapsed during the assembly process (Peona et al. 2018). Assembly contiguity and completeness (as
23 judged by BUSCO scores) of the *O. melanoleuca* assembly compared favourably to other songbird
24 genome assemblies (**Tab. 2**).

25 **Mitogenome assembly**

26 MitoFinder and MitoVGP assembled mitogenomes of 16,944 bp and 18,631 bp length, respectively.
27 The mitochondrial contigs assembled by the two pipelines were congruent, except for 9 single base
28 pair mismatches, for a 1,827 bp long insert in the MitoVGP assembly and of a 141 bp long insert in the
29 MitoFinder assembly. We decided to not consider either of these inserts in the final mitogenome
30 assembly for the following reasons. First, neither of the inserts was observed in the mitogenomes of
31 isabelline and northern wheatear. For the long insert in the MitoVGP assembly, moreover, the
32 coverage of short reads mapped to the MitoVGP assembly was strongly reduced (**Fig. S1**), and the
33 insertion constituted a partial duplication of *nd6*, duplications of two tRNAs (Glu, Pro) and a partial
34 duplication of the control region likely caused by an assembly artefact. The short insert in the
35 MitoFinder assembly was not observed in the other wheatear mitogenomes, and if real, we would
36 expect long reads to cover this insert. Because base calling based on short reads is expected to have
37 higher quality, we retained the MitoFinder assembly, but without the 141 bp insert as final
38 mitogenome.

39 The final mitogenome (as also both original assemblies) contained all 13 protein-coding
40 genes, two rRNAs, and 22 tRNAs (**Fig. 2**). All genes, except eight tRNAs and *nd6*, were located on the

1 heavy DNA strand. Both gene order and strandedness were concordant with those observed in
2 northern wheatear (*O. oenanthe*) (Wang et al. 2020).

3 Repetitive element annotation

4 The *de novo* identification of repetitive elements resulted in the characterisation of 572 raw
5 consensus sequences from RepeatModeler2 and 16 satellite DNA consensus sequences from
6 RepeatExplorer2. The consensus sequences from RepeatModeler2 were filtered from tandem repeats
7 and protein-coding genes. This resulted in a final library of 477 consensus sequences (**File S1**).
8 Among these consensus sequences, RepeatModeler2 classified 226 sequences as LTR
9 retrotransposons, 98 as LINE retrotransposons, 21 as DNA transposons, 5 as SINE retrotransposons,
10 and 112 sequences were unclassified (“unknown”).

11 The genome assembly annotation run with RepeatMasker using the repeat library produced
12 here and merged with already known avian repeats showed that ~10 % of the assembled genome is
13 repetitive (**Fig. 3A, Tab. S1, File S2**). This finding indicates that many repeats collapsed during the
14 genome assembly process. An example of this were satDNAs that represented ~0.8 % of the
15 sequenced reads but only < 0.3 % of the genome assembly, suggesting that satDNA repeats (such as
16 in (peri-)centromeric and (sub-)telomeric regions) are the most collapsed repeats. Most of the
17 repeats annotated were LTR and LINE retrotransposons (**Fig. 3A**). While it is common to find LINEs
18 as most abundant TEs in avian genomes (Galbraith et al. 2021; Kapusta and Suh 2017; Manthey et al.
19 2018; Peona, Blom et al. 2021), it is less common to find so similar percentages of LINE and LTR
20 retrotransposons. This is especially true for a male genome assembly such as the present one here
21 that does not include the W chromosome which is highly enriched in LTRs and acts as a refugium for
22 most of the full-length genomic LTR elements in birds (Peona et al. 2021; Warmuth et al. 2022). The
23 transposable element landscape (**Fig. 3B**) suggests that LINE retrotransposons experienced a drop in
24 their genomic accumulation in recent times (0-5 % divergence; **Fig. 3B**), whereas LTR
25 retrotransposons kept accumulating at the same rate. Such a recent replacement of LINE
26 retrotransposon activity with a diversity of LTR retrotransposons has been noted in other songbirds
27 and seems to have occurred independently in the so far analysed passerine families, i.e., estrildid
28 finches (Warren et al. 2010, Boman et al. 2019), flycatchers (Suh et al. 2018), crows (Weissensteiner
29 et al. 2020), and birds-of-paradise (Peona et al. 2021). Finally, the satDNA landscape (**Fig. 3B**) shows
30 that satDNA arrays experienced differential amplification in copies number in recent times (0-10 %
31 divergence), implying fast evolution of this genomic fraction in the genome (Peona et al. 2022).

32

1 Transcriptome sequencing, genome annotation, and gene function prediction

2 Iso-Seq sequencing yielded a total of 4,627,382 CCS reads (125,633-1,087,892 reads per tissue, **Tab.**
3 **3**). This resulted in numbers of high-quality isoforms ranging from 16,078 to 80,600 per tissue. On
4 average 8'833 genes were predicted per tissue, ranging from 4,772 in muscle to 10,924 in liver.
5 Transcriptome completeness evaluated through BUSCO ranged from 31.2 % to 57.5 % complete
6 BUSCO per tissue (**Tab. 3**).

7 The Iso-Seq transcriptomes were then used as splice evidence in GeMoMa to perform a
8 predominantly homology-based annotation of the reference genome. We predicted 18,143 protein-
9 coding genes with a total of 320,754 exons and 289,421 introns. The number of exons, CDS, and
10 introns was higher for our *O. melanoleuca* annotation compared to the annotations of other songbirds,
11 such as *Junco hyemalis*, *Fringilla coelebs*, *Melospiza melodia*, *Taeniopygia guttata*, *Ficedula albicollis*,
12 *Manacus vitellinus*, and *Geospiza fortis* (**Tab. 2**). Mean gene length, CDS length, exon length, and
13 number of exons per gene, on the other hand, were in the range of values obtained for the
14 abovementioned songbird annotations (**Tab. 3**). 17,'553 (96.7 %) of the 18,'143 predicted genes were
15 annotated with protein families or function assignment. 12,'472 (68.7 %) genes obtained a GO term
16 assignment through InterProScan. The most abundant GO terms were associated with “cell part”,
17 “cell” and “membrane” in the cellular component category, “binding” in the molecular function
18 category and “cellular metabolic process” or “metabolic process” in the biological process category
19 (**Fig. S2**). BUSCO completeness of the final annotation as judged from avian BUSCO (n=8,338) was
20 98.0 %, with 97.4 % single copy BUSCO, 0.6 % duplicated BUSCO, 0.6 % fragmented BUSCO, and 1.5
21 % missing BUSCO. This suggests an accurate and rather complete annotation.

22 Data Availability

23 All data, including the assembly, its annotation, and the original sequencing data are available on the
24 European Nucleotide Archive under project accession PRJNA937434. Code for the repeat analysis is
25 available on <https://github.com/ValentinaBoP/WheatearGenomeAnalysis>.

26 Supplemental Material is available at figshare: <https://doi.org/10.25387/g3.22209697>.

27 Acknowledgements

28 We warmly thank Marta Burri for preparing RNA, and Giulio Formenti, Remi Allio, and Lauren
29 Coombe for support with computational questions. We are indebted to NGI Uppsala, namely Mai-Britt
30 Mosbech and Olga Vinnere Pettersson for UHMW DNA extraction and long-read sequencing and Ignas
31 Bunikis for running the primary assembly, as well as NGI Stockholm for the preparation of linked-
32 read and Hi-C sequencing data. Finally, we thank the NGS platform at the University of Berne, namely
33 Pamela Nicholson and Catia Coito for the preparation of Iso-Seq data. Computations were performed
34 on resources provided by the Swedish National Infrastructure for Computing (SNIC) through the
35 Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) and the High-
36 Performance Computing Cluster EVE, a joint effort of the Helmholtz Centre for Environmental
37 Research (UFZ) and the German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-

1 Leipzig. We thank the administration and support staff of EVE: Thomas Schnicke and Ben Langenberg
2 (UFZ), and Christian Krause (iDiv).

3 **Conflict of Interest**

4 The authors declare no conflict of interest.

5 **Funder Information**

6 The present project was supported by German Research Foundation (DFG), grant number BU3456/3-
7 1 to RB and the National Research Fund (FNR) Luxembourg, grant number 14575729, to DL; OMPG
8 was supported by the Swedish Research Council Vetenskapsrådet (grant number 2020-03866); VP
9 was supported via grants to AS from the Swedish Research Council Vetenskapsrådet (grant number
10 2020-04436) and the Swedish Research Council Formas (2017-01597). NAK was supported by a
11 Georg Foster Research Stipend of the Alexander von Humboldt Foundation. Part of the analyses were
12 enabled by resources provided by the Swedish National Infrastructure for Computing (SNIC) at
13 Uppsala partially funded by the Swedish Research Council through grant agreement no. 2018-05973.
14

1 Literature cited

- 2 Alaei Kakhki N, Schweizer M, Lutgen D, Bowie RCK, Shirihai H, Suh A, Schielzeth H, Burri R. In press.
3 A phylogenomic assessment of processes underpinning convergent evolution in open-habitat
4 chats. *Mol. Biol. Evol* (in press, bioRxiv doi:10.1101/2022.06.21.496980).
- 5 Aliabadian M, Kaboli M, Förchler MI, Nijman V, Chamani A, Tillier A, Prodon R, Pasquet E, Ericson
6 PGP, Zuccon D. 2012. Convergent evolution of morphological and ecological traits in the open-
7 habitat chat complex (Aves, Muscicapidae: Saxicolinae). *Molecular Phylogenetics and Evolution*.
8 65(2012):35–45.
- 9 Allio R, Schomaker-Bastos A, Romiguier J, Prosdocimi F, Nabholz B, Delsuc F. 2020. MitoFinder:
10 Efficient automated large-scale extraction of mitogenomic data in target enrichment
11 phylogenomics. *Molecular ecology resources*. 20(4):892–905. eng. doi:10.1111/1755-
12 0998.13160.
- 13 Babraham Bioinformatics: Cambridge. 2012. FastQC; Version 0.10.1: A Quality Control Tool for High
14 throughput Sequence Data; Babraham Bioinformatics. Cambridge, UK.
- 15 Bao W, Kojima KK, Kohany O. 2015. Repbase Update, a database of repetitive elements in eukaryotic
16 genomes. *Mobile DNA*. 6:11. eng. doi:10.1186/s13100-015-0041-9.
- 17 Bellott DW, Skaletsky H, Cho T-J, Brown L, Locke D, Chen N, Galkina S, Pyntikova T, Koutseva N,
18 Graves T, et al. 2017. Avian W and mammalian Y chromosomes convergently retained dosage-
19 sensitive regulators. *Nature genetics*. 49(3):387–394. eng. doi:10.1038/ng.3778.
- 20 Boman J, Frankl-Vilches C, da Silva dos Santos, Michelly, de Oliveira, Edivaldo H. C., Gahr M, Suh A.
21 2019. The Genome of Blue-Capped Cordon-Bleu Uncovers Hidden Diversity of LTR
22 Retrotransposons in Zebra Finch. *Genes*. 10(4):301.
- 23 Boman J, Frankl-Vilches C, da Silva dos Santos, Michelly, Oliveira EHC de, Gahr M, Suh A, da Silva dos
24 Santos, Michelly. 2019. The Genome of Blue-Capped Cordon-Bleu Uncovers Hidden Diversity of
25 LTR Retrotransposons in Zebra Finch. *Genes*. 10(4). eng. doi:10.3390/genes10040301.
- 26 Buerkle CA, Lexer C. 2008. Admixture as the basis for genetic mapping. *Trends in Ecology &*
27 *Evolution*. 23(12):686–694. doi:10.1016/j.tree.2008.07.008.
- 28 Cabanettes F, Klopp C. 2018. D-GENIES: Dot plot large genomes in an interactive, efficient and
29 simple way. *PeerJ*. 6:e4958. eng. doi:10.7717/peerj.4958.
- 30 Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: An ultra-fast all-in-one FASTQ preprocessor.
31 *Bioinformatics*. 34(17):i884–i890. eng. doi:10.1093/bioinformatics/bty560.
- 32 Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-
33 Balderas R, Morales-Cruz A, et al. 2016. Phased diploid genome assembly with single-molecule
34 real-time sequencing. *Nat Meth*. 13(12):1050–1054. doi:10.1038/nmeth.4035.
- 35 Cornetti L, Valente LM, Dunning LT, Quan X, Black RA, Hébert O, Savolainen V. 2015. The Genome of
36 the "Great Speciator" Provides Insights into Bird Diversification. *Genome Biol Evol*. 7(9):2680–
37 2691. eng. doi:10.1093/gbe/evv168.
- 38 Dainat J. 2019. AGAT: Another Gff Analysis Toolkit to handle annotations in any GTF/GFF format.
39 (Version v0.7.0). Zenodo. doi:10.5281/zenodo.3552717.

- 1 Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, Shamim MS, Machol I, Lander
2 ES, Aiden AP, et al. 2017. De novo assembly of the *Aedes aegypti* genome using Hi-C yields
3 chromosome-length scaffolds. *Science (New York, N.Y.)*. 356(6333):92–95. eng.
4 doi:10.1126/science.aal3327.
- 5 Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, Aiden EL. 2016. Juicer Provides
6 a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell systems*. 3(1):95–98.
7 eng. doi:10.1016/j.cels.2016.07.002.
- 8 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
9 *Nucl. Acids Res.* 32(5):1792–1797. doi:10.1093/nar/gkh340.
- 10 Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, Kunstner A, Makinen H,
11 Nadachowska-Brzyska K, Qvarnstrom A, et al. 2012. The genomic landscape of species
12 divergence in *Ficedula* flycatchers. *Nature*. 491(7426):756–760. doi:10.1038/nature11584.
- 13 Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. 2020. RepeatModeler2 for
14 automated genomic discovery of transposable element families. *Proceedings of the National
15 Academy of Sciences of the United States of America*. 117(17):9451–9457. eng.
16 doi:10.1073/pnas.1921046117.
- 17 Formenti G, Rhie A, Balacco J, Haase B, Mountcastle J, Fedrigo O, Brown S, Capodiferro MR, Al-Ajli
18 FO, Ambrosini R, et al. 2021. Complete vertebrate mitogenomes reveal widespread repeats and
19 gene duplications. *Genome Biol.* 22(1):120. eng. doi:10.1186/s13059-021-02336-9.
- 20 Friis G, Vizueta J, Ketterson ED, Milá B. 2022. A high-quality genome assembly and annotation of the
21 dark-eyed junco *Junco hyemalis*, a recently diversified songbird. *G3 (Bethesda)*. 12(6). eng.
22 doi:10.1093/g3journal/jkac083.
- 23 Galbraith JD, Kortschak RD, Suh A, Adelson DL. 2021. Genome Stability Is in the Eye of the Beholder:
24 CR1 Retrotransposon Activity Varies Significantly across Avian Diversity. *Genome Biol Evol.*
25 13(12). eng. doi:10.1093/gbe/evab259.
- 26 Gompert Z, Mandeville EG, Buerkle CA. 2017. Analysis of Population Genomic Data from Hybrid
27 Zones. *Annual Review of Ecology, Evolution, and Systematics*; [accessed 2017 Aug 21]. 48:207–
28 229. doi:10.1146/annurev-ecolsys-110316-022652.
- 29 Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. 2020. Identifying and removing
30 haplotypic duplication in primary genome assemblies. *Bioinformatics*. 36(9):2896–2898. eng.
31 doi:10.1093/bioinformatics/btaa025.
- 32 Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUASt: Quality assessment tool for genome
33 assemblies. *Bioinformatics*. 29(8):1072–1075. eng. doi:10.1093/bioinformatics/btt086.
- 34 Haffer J. 1977. Secondary contact zones of birds in Northern Iran. *Bonner Zoologische
35 Monographien*. (10):1–64.
- 36 International Chicken Genome Sequencing Consortium. 2004. Sequence and comparative analysis of
37 the chicken genome provide unique perspectives on vertebrate evolution. *Nature*.
38 432(7018):695–716. eng. doi:10.1038/nature03154.
- 39 Jackman SD, Coombe L, Chu J, Warren RL, Vandervalk BP, Yeo S, Xue Z, Mohamadi H, Bohlmann J,
40 Jones SJM, et al. 2018. Tigmint: Correcting assembly errors using linked reads from large
41 molecules. *BMC bioinformatics*. 19(1):393. eng. doi:10.1038/nmeth.4035.

- 1 Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G,
2 et al. 2014. InterProScan 5: Genome-scale protein function classification. *Bioinformatics*.
3 30(9):1236–1240. eng. doi:10.1093/bioinformatics/btu031.
- 4 Kapusta A, Suh A. 2017. Evolution of bird genomes—a transposon's-eye view. *Annals of the New
5 York Academy of Sciences*. 1389(1):164–185. doi:10.1111/nyas.13295.
- 6 Kapusta A, Suh A, Feschotte C. 2017. Dynamics of genome size evolution in birds and mammals.
7 *Proceedings of the National Academy of Sciences of the United States of America*. 114(8):E1460-
8 E1469.
- 9 Kawakami T, Backström N, Burri R, Husby A, Olason P, Rice AM, Ålund M, Qvarnström A, Ellegren H.
10 2014. Estimation of linkage disequilibrium and interspecific gene flow in *Ficedula* flycatchers by
11 a newly developed 50k single-nucleotide polymorphism array. *Molecular Ecology Resources*.
12 14(6):1248–1260. doi:10.1111/1755-0998.12270.
- 13 Kawakami T, Smeds L, Backström N, Husby A, Qvarnström A, Mugal CF, Olason P, Ellegren H. 2014. A
14 high-density linkage map enables a second-generation collared flycatcher genome assembly and
15 reveals the patterns of avian recombination rate variation and chromosomal evolution.
16 *Molecular Ecology*. 23(16):4035–4058. doi:10.1111/mec.12810.
- 17 Keilwagen J, Hartung F, Grau J. 2019. GeMoMa: Homology-Based Gene Prediction Utilizing Intron
18 Position Conservation and RNA-seq Data. *Methods Mol Biol*. 1962:161–177. eng.
19 doi:10.1007/978-1-4939-9173-0_9.
- 20 Keilwagen J, Hartung F, Paulini M, Twardziok SO, Grau J. 2018. Combining RNA-seq data and
21 homology-based gene prediction for plants, animals and fungi. *BMC bioinformatics*. 19(1):189.
22 eng. doi:10.1186/s12859-018-2203-5.
- 23 Kuo RI, Cheng Y, Zhang R, Brown JWS, Smith J, Archibald AL, Burt DW. 2020. Illuminating the dark
24 side of the human transcriptome with long read transcript sequencing. *BMC genomics*.
25 21(1):751. eng. doi:10.1186/s12864-020-07123-7.
- 26 Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*. 34(18):3094–
27 3100. eng. doi:10.1093/bioinformatics/bty191.
- 28 Li H. 2021. New strategies to improve minimap2 alignment accuracy. *Bioinformatics*. 37(23):4572–
29 4574. eng. doi:10.1093/bioinformatics/btab705.
- 30 Lutgen D, Ritter R, Olsen R-A, Schielzeth H, Gruselius J, Ewels P, García JT, Shirihai H, Schweizer M,
31 Suh A, et al. 2020. Linked-read sequencing enables haplotype-resolved resequencing at
32 population scale. *Molecular Ecology Resources*. 20(5):1311–1322. doi:10.1111/1755-
33 0998.13192.
- 34 Manthey JD, Moyle RG, Boissinot S. 2018. Multiple and Independent Phases of Transposable Element
35 Amplification in the Genomes of Piciformes (Woodpeckers and Allies). *Genome Biol Evol*.
36 10(6):1445–1456. eng. doi:10.1093/gbe/evy105.
- 37 Marçais G, Kingsford C. 2011. A fast, lock-free approach for efficient parallel counting of occurrences
38 of k-mers. *Bioinformatics*. 27(6):764–770. eng. doi:10.1093/bioinformatics/btr011.
- 39 Novák P, Neumann P, Macas J. 2020. Global analysis of repetitive DNA from unassembled sequence
40 reads using RepeatExplorer2. *Nature Protocols*. 15(11):3745–3776. eng. doi:10.1038/s41596-
41 020-0400-y.

- 1 Panov EN. 2005. Wheaters of the Palearctic. Ecology, Behaviour and Evolution of the genus
2 *Oenanthe*. Sofia-Moscow: Pensoft (Wilson MG, editor. Series Faunistica).
- 3 Paysan-Lafosse T, Blum M, Chuguransky S, Grego T, Pinto BL, Salazar GA, Bileschi ML, Bork P, Bridge
4 A, Colwell L, et al. 2022. InterPro in 2022. Nucleic acids research. eng. doi:10.1093/nar/gkac993.
- 5 Peona V, Blom MPK, Xu L, Burri R, Sullivan S, Bunikis I, Liachko I, Haryoko T, Jønsson KA, Zhou Q, et
6 al. 2021. Identifying the causes and consequences of assembly gaps using a multiplatform
7 genome assembly of a bird-of-paradise. Molecular ecology resources. 21(1):263–286. eng.
8 doi:10.1111/1755-0998.13252.
- 9 Peona V, Kutschera VE, Blom MPK, Irestedt M, Suh A. 2022. Satellite DNA evolution in Corvoidea
10 inferred from short and long reads. Molecular Ecology. eng. doi:10.1111/mec.16484.
- 11 Peona V, Palacios-Gimenez OM, Blommaert J, Liu J, Haryoko T, Jønsson KA, Irestedt M, Zhou Q, Jern
12 P, Suh A. 2021. The avian W chromosome is a refugium for endogenous retroviruses with likely
13 effects on female-biased mutational load and genetic incompatibilities. Philosophical
14 transactions of the Royal Society of London. Series B, Biological sciences. 376(1833):20200186.
15 eng. doi:10.1098/rstb.2020.0186.
- 16 Peona V, Weissensteiner MH, Suh A. 2018. How complete are “complete” genome assemblies?: —An
17 avian perspective. Molecular Ecology Resources. 18(6):1188–1195. doi:10.1111/1755-
18 0998.12933.
- 19 Rhie A, McCarthy SA, Fedrigo O, Damas J, Formenti G, Koren S, Uliano-Silva M, Chow W,
20 Functamman A, Kim J, et al. 2021. Towards complete and error-free genome assemblies of all
21 vertebrate species. Nature. 592(7856):737–746. eng. doi:10.1038/s41586-021-03451-0.
- 22 Robinson JT, Turner D, Durand NC, Thorvaldsdóttir H, Mesirov JP, Aiden EL. 2018. Juicebox.js
23 Provides a Cloud-Based Visualization System for Hi-C Data. Cell systems. 6(2):256-258.e1. eng.
24 doi:10.1016/j.cels.2018.01.001.
- 25 Schweizer M, Warmuth V, Alaei Kakhki N, Aliabadian M, Förschler MI, Shirihai H, Suh A, Burri R.
26 2019. Parallel plumage color evolution and pervasive hybridization in wheatears. Journal of
27 Evolutionary Biology. 32(1):100–110.
- 28 Schweizer M, Warmuth VM, Alaei Kakhki N, Aliabadian M, Förschler M, Shirihai H, Ewels P, Gruselius
29 J, Olsen R-A, Schielzeth H, et al. 2019. Genome-wide evidence supports mitochondrial
30 relationships and pervasive parallel phenotypic evolution in open-habitat chats. Molecular
31 Phylogenetics and Evolution. 139:106568.
- 32 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing
33 genome assembly and annotation completeness with single-copy orthologs. Bioinformatics;
34 [accessed 12/25/2019]. 31(19):3210–3212. doi:10.1093/bioinformatics/btv351.
- 35 Smit AFA, Hubley R, Green P. 1996-2010. RepeatMasker Open 3.3.0.
- 36 Stecher G, Tamura K, Kumar S. 2020. Molecular Evolutionary Genetics Analysis (MEGA) for macOS.
37 Mol Biol Evol. 37(4):1237–1239. eng. doi:10.1093/molbev/msz312.
- 38 Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. 2021. The Dfam community resource of
39 transposable element families, sequence models, and genome annotations. Mobile DNA. 12(1):2.
40 eng. doi:10.1186/s13100-020-00230-y.

- 1 Stothard P, Wishart DS. 2005. Circular genome visualization and exploration using CGView.
2 Bioinformatics. 21(4):537–539. eng. doi:10.1093/bioinformatics/bti054.
- 3 Suh A, Smeds L, Ellegren H. 2018. Abundant recent activity of retrovirus-like retrotransposons
4 within and among flycatcher species implies a rich source of structural variation in songbird
5 genomes. *Molecular Ecology*. in press. doi:10.1111/mec.14439.
- 6 Uniprot Consortium. 2019. UniProt: A worldwide hub of protein knowledge. *Nucleic acids research*.
7 47(D1):D506–D515. eng. doi:10.1093/nar/gky1049.
- 8 Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC. 2017.
9 GenomeScope: Fast reference-free genome profiling from short reads. *Bioinformatics*.
10 33(14):2202–2204. eng. doi:10.1093/bioinformatics/btx153.
- 11 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J,
12 Young SK, et al. 2014. Pilon: An integrated tool for comprehensive microbial variant detection
13 and genome assembly improvement. *PLoS ONE*. 9(11):e112963. eng.
14 doi:10.1371/journal.pone.0112963.
- 15 Wang E, Zhang D, Braun MS, Hotz-Wagenblatt A, Pärt T, Arlt D, Schmaljohann H, Bairlein F, Lei F,
16 Wink M. 2020. Can Mitogenomes of the Northern Wheatear (*Oenanthe oenanthe*) Reconstruct
17 Its Phylogeography and Reveal the Origin of Migrant Birds? *Sci Rep*. 10(1):9290. eng.
18 doi:10.1038/s41598-020-66287-0.
- 19 Warmuth VM, Weissensteiner MH, Wolf JBW. 2022. Accumulation and ineffective silencing of
20 transposable elements on an avian W Chromosome. *Genome research*. 32(4):671–681. eng.
21 doi:10.1101/gr.275465.121.
- 22 Warren RL, Yang C, Vandervalk BP, Behsaz B, Lagman A, Jones SJM, Birol I. 2015. LINKS: Scalable,
23 alignment-free scaffolding of draft genomes with long reads. *GigaScience*. 4:35. eng.
24 doi:10.1186/s13742-015-0076-3.
- 25 Warren WC, Clayton DF, Ellegren H, Arnold AP, Hillier LW, Kunster A, Searle S, White S, Vilella AJ,
26 Fairly S, et al. 2010. The genome of a songbird. *Nature*. 464(7289):757–762.
- 27 Weissensteiner MH, Bunikis I, Catalán A, Francoijs K-J, Knief U, Heim W, Peona V, Pophaly SD,
28 Sedlazeck FJ, Suh A, et al. 2020. Discovery and population genomics of structural variation in a
29 songbird genus. *Nat Commun*. 11(1):3403. eng. doi:10.1038/s41467-020-17195-4.
- 30 Yeo S, Coombe L, Warren RL, Chu J, Birol I. 2018. ARCS: Scaffolding genome drafts with linked reads.
31 *Bioinformatics*. 34(5):725–731. eng. doi:10.1093/bioinformatics/btx675.

32
33

Table 1. Assembly statistics for different versions of the *O. melanoleuca* genome.

		Falcon unzip, Arrow	+ Pilon, purge_dups	+ Tigmint	+ 3D DNA (all)	+ 3D DNA (chrom)
Basic stats	No. contigs/scaffolds*	1,681	381	383	588*	32*
	No. contigs/scaffolds* > 50 kb	1,610	347	348	143*	31*
	Assembly length (Gb)	1.29	1.04	1.04	1.04*	1.00*
	Contig/scaffold* N50 (Mb)	8.6	13.5	12.6	69.6*	69.7*
	Contig/scaffold* L50	35	23	24	6*	5*
	Largest contig/scaffold* (Mb)	45.3	45.3	45.3	148.4*	148.4*
BUSCO	Complete (%)	96.9	96.4	96.4	96.2	95.5
	Complete single-copy (%)	90.6	95.9	95.9	95.7	95.1
	Complete duplicated (%)	6.3	0.5	0.5	0.5	0.4
	Fragmented (%)	0.7	0.7	0.7	0.9	0.9
	Missing (%)	2.4	2.9	2.9	2.9	3.6

* Where numbers concern scaffolds instead of contigs, this is indicated by an asterisk.

Table 2. Comparison of genome assembly and annotation summary statistics of *Oenanthe melanoleuca* with other songbird species (*Junco hyemalis*, *Fringilla coelebs*, *Melospiza melodia*, *Taeniopygia guttata*, *Ficedula albicollis*, *Manacus vitellinus*, and *Geospiza fortis*). Modified from Friis et al. (2022).

	<i>Oenanthe</i>	<i>Junco</i>	<i>Fringilla</i>	<i>Melospiza</i>	<i>Taeniopygia</i>	<i>Ficedula</i>	<i>Manacus</i>	<i>Geospiza</i>	
Genome assembly length (Gb)	1.04	0.99	0.99	1.36	1.22	1.1	1.17	1.04	
Genome contig N50 (kb)	7,700	75	67	8,300	38	410	194	30	
Genome BUSCO scores (%)	C	95.5	95.4	94.1	87.9	93.8	96.5	96.1	96.0
	S	95.1	95.2	93.8	87.3	91.9	96	94.6	95.6
	D	0.4	0.2	0.3	0.6	1.9	0.5	1.5	0.4
	F	0.9	1.6	2.0	7.2	2.3	0.8	1	1.2
	M	3.6	3.0	4.0	5.0	3.9	2.7	2.9	2.8
No. of genes	18,143	19,026	17,703	15,086	17,561	16,763	18,976	14,399	
Mean gene length (bp)	28,23218	15,402	15,818	14,457	26,458	31,394	27,847	30,164	
No. of CDS	31,333	23,245	17,703	15,086	17,561	16,763	18,976	14,399	
Mean CDS length (bp)	1682	1,647	1,679	1,325	1,677	1,942	1,929	1,766	
No. of exons	320,754	229,210	221,872	131,940	171,767	189,043	190,390	164,721	
Mean exon length (bp)	164	167	165	153	255	253	264	195	
Mean no. exons/gene	102	9.9	10.2	8.7	10.3	12.2	11.5	11.4	
No. of introns	289,421	205,965	200,041	116,724	153,909	171,236	171,089	149,563	

BUSCO parameters are C: complete genes; S: complete and single-copy genes; D: complete and duplicated genes; F: fragmented genes; M: missing genes.

1

Table 3. Iso-Seq data characterization and transcriptome completeness.

		Brain	Heart	Kidney	Liver	Lung	Muscle	Testis
Transcriptome	No. of CCS reads	847,617	253,468	723,158	1,087,892	1,061,936	125,633	527,678
	High-quality isoforms	73,422	80,600	45,097	47,491	28,508	16,078	44,605
	Low-quality isoforms	734	844	616	384	151	94	284
	No. of genes	10,449	10,448	9,063	10,924	6,564	4,772	9,613
	Mean gene length (bp)	24,193	20,119	16,350	15,125	18,528	17,397	17,415
	No. of CDS	27,449	28,747	25,790	27,202	13,551	8,447	23,009
	Mean CDS length (bp)	972	985	932	823	894	980	960
	No. of exons	231,169	222,791	235,989	194,325	108,084	69,859	184,794
	Mean exon length (bp)	246	248	223	225	221	224	209
	Mean no. of exons/mRNA	8.4	8.2	7.9	7.1	8.0	8.3	8.0
BUSCO	Complete (%)	56.80	57.50	48.30	49.40	38.30	31.20	49.3
	Single-copy (%)	40.30	39.50	33.60	34.70	31.1	27.00	34.6
	Duplicated (%)	16.50	18.00	14.70	14.70	7.20	4.20	14.70
	Fragmented (%)	2.90	2.10	2.60	3.20	2.00	1.10	2.30
	Missing (%)	40.30	40.60	49.10	47.40	59.70	67.70	48.40

2

3

4

1 **Figure Legends**

2 **Figure 1.** Eastern black-eared wheatear (*Oenanthe melanoleuca*). The species sports a white-throated
3 (left; Agii Pantes, Greece, June 2022) and a black-throated phenotype (right; Lesvos, Greece, May
4 2017) in males. © Reto Burri

5
6 **Figure 2.** Circular sketch map of the *O. melanoleuca* mitogenome assembly. The outer circle shows
7 coding sequences (purple), rRNAs (pink), and tRNAs (red). The black trace on the middle circle
8 indicates GC content. On the inner circle, positive and negative GC skews in nucleotide composition
9 are indicated by green and magenta, respectively.

10
11 **Figure 3.** Repeat annotation landscapes. **A)** Pie-chart summarizing the transposable element content
12 annotated in the genome assembly. **B)** Transposable element landscape. The divergence between
13 interspersed repeat copies and their consensus sequences is shown on the X-axis as genetic distance
14 calculated using the Kimura 2-parameter distance. The percentage of the genome assembly occupied
15 by transposable elements is shown on the Y-axis. **C)** Satellite DNA landscape. The divergence between
16 the satellite DNA consensus sequences and sequences annotated in the short-read library is shown
17 on the X-axis as genetic distance calculated using the Kimura 2-parameter distance. The percentage
18 of the genome (short reads) annotated as satellite DNA is shown on the Y-axis.

19
20



Figure 1
160x53 mm (x DPI)

1
2
3
4

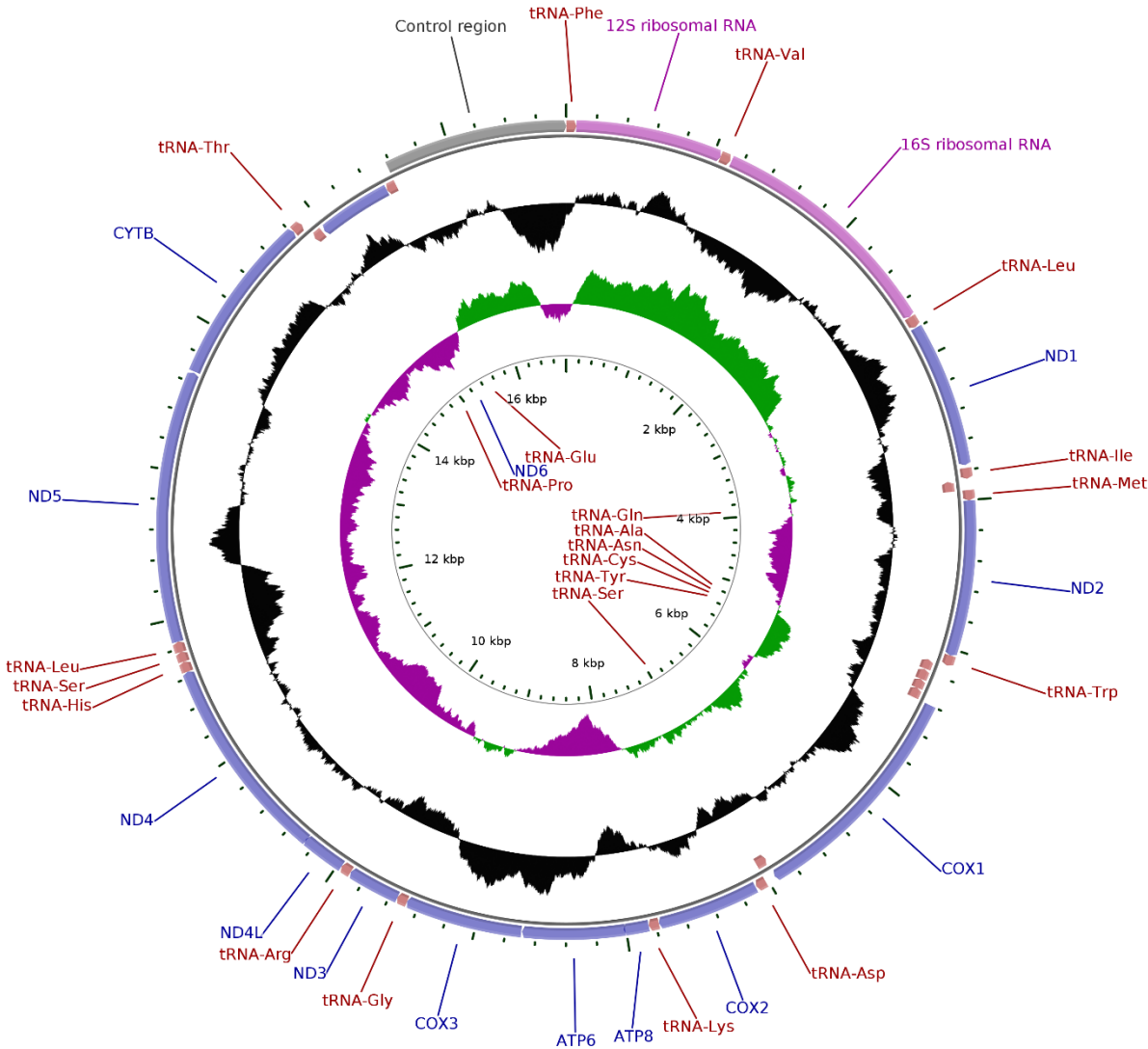
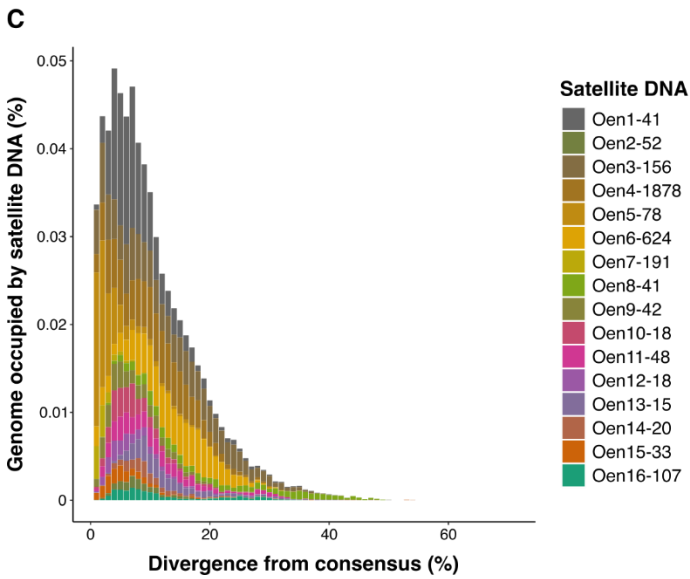
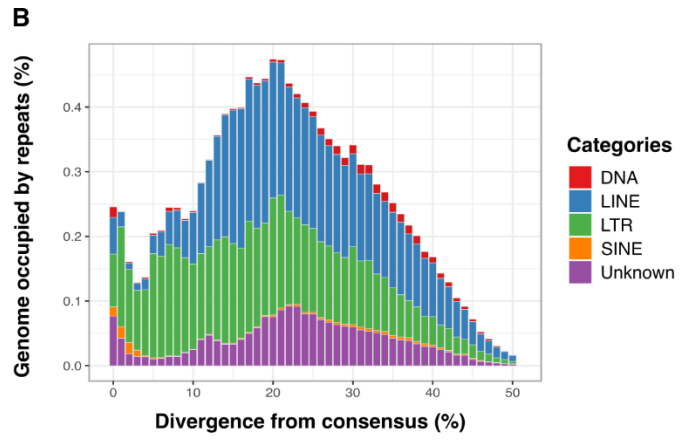
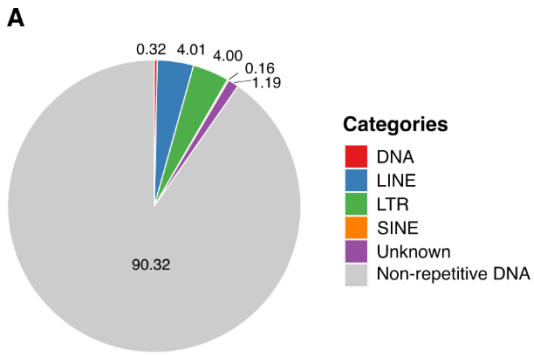


Figure 2
160x145 mm (x DPI)

1
2
3
4



1
2
3

Figure 3
160x135 mm (x DPI)